

1 **Title:**

2 **Greengenes: Chimera-checked 16S rRNA gene database and workbench compatible**
3 **with ARB**

4 **Authors:**

5 DeSantis, T.Z.¹

6 Hugenholtz, P.²

7 Larsen, N.³

8 Rojas, M.⁴

9 Brodie, E.L.¹

10 Keller, K.⁵

11 Huber, T.⁶

12 Dalevi, D.⁷

13 Hu, P.¹

14 Andersen, G.L.¹

15

16 ¹Center for Environmental Biotechnology

17 Lawrence Berkeley National Laboratory

18 1 Cyclotron Road, Mail Stop 70A-3317

19 Berkeley, CA 94720

20 USA

21 ²Microbial Ecology Program

22 DOE Joint Genome Institute

23 2800 Mitchell Drive Bldg 400-404

24 Walnut Creek, CA 94598

25 USA

26 ³Danish Genome Institute

27 Gustav Wieds vej 10 C

28 DK-8000 Aarhus C

29 Denmark

30 ⁴Department of Bioinformatics

31 Baylor University

32 P.O. Box 97356, 1311 S. 5th St.

33 Waco, TX 76798-7356

34 USA

35 ⁵Department of Bioengineering

36 University of California

37 Berkeley, CA 94720

38 USA

39 ⁶Departments of Biochemistry and Mathematics

40 The University of Queensland

41 Brisbane Qld 4072

42 Australia

43 ⁷Department of Computer Science

44 Chalmers University of Technology

1 SE-412 96
2 Göteborg, Sweden

3
4 To whom correspondence should be addressed:
5 GLAndersen@lbl.gov

6
7
8 Abstract

9 A 16S rRNA gene database (<http://greengenes.lbl.gov>) addresses limitations of public
10 repositories by providing chimera-screening, standard alignments and taxonomic
11 classification using multiple published taxonomies. It was revealed that incongruent
12 taxonomic nomenclature exists among curators even at the phylum-level. Putative
13 chimeras were identified in 3% of environmental sequences and 0.2% of records derived
14 from isolates. Environmental sequences were classified into 100 phylum-level lineages
15 within the Archaea and Bacteria.

1

2 Comparative analysis of 16S small sub-unit rRNA genes (hereafter abbreviated as
3 16S) is commonly applied to survey the constituents of microbial communities (4, 13, 23,
4 24), to infer bacterial and archaeal evolution (14, 19), and to design monitoring and
5 analysis tools such as microarrays (5, 10, 17, 20, 29, 30). Because the production rate of
6 16S sequence records from uncultured organisms now exceeds that from their cultured
7 counterparts, taxonomic placement of sequences lags behind. In fact, 43% of full-length
8 16S records in GenBank are amalgamated into pseudo-divisions “environmental
9 samples” or “unclassified”. Annotation styles are inconsistent creating barriers for
10 computational categorization of biological sources. Furthermore, since rRNA genes from
11 environmental DNA are usually PCR amplified, it is suspected that many clandestine
12 chimeric sequences are intercalated into the public databases. From a small sample of
13 1,399 sequence records from known phyla, it was estimated that 3% of the public data
14 might contain chimeras (2). The permeation of these poor quality data along with barriers
15 in exchanging nomenclature have led to several conflicting taxonomies. The probability
16 of mistakenly adopting a chimeric sequence in a phylogenetic inference or as a reference
17 for probe/primer design is rising noticeably. Lastly, ARB (21) database administration
18 needs to be streamlined for those who maintain 16S collections on their local computers.

19 Greengenes addresses these concerns by providing four features: a standardized
20 set of descriptive fields, taxonomic assignment, chimera screening and ARB
21 compatibility. Heuristics are used to consider the author’s annotations and categorize
22 each source as a named or unnamed isolate, an unnamed symbiont, or an uncultured
23 organism. Other standard descriptors include sequence quality measurements, authors,

1 and a “study_id” that links all the records associated with a project. Greengenes
2 maintains a consistent multiple sequence alignment of both archaeal and bacterial 16S
3 genes to facilitate taxonomic placement. Taxonomy proposed by independent curators
4 (NCBI, RDP (Bergey’s) (7), Wolfgang Ludwig (21), Phil Hugenholtz (16) and Norm
5 Pace (23)) is tracked to promote user awareness of several estimations of phylogenetic
6 descent allowing a balanced approach to node nomenclature when generating
7 dendrograms. Comprehensive chimera assessment is a distinguishing characteristic of the
8 Greengenes data assembly process. Each sequence is scored for chimeric potential, a
9 break point is estimated and parent sequences are identified. Furthermore, since
10 biologists often collect and visualize 16S relationships using the freely available ARB
11 software, Greengenes simplifies the chore of keeping a research group’s private ARB
12 database current by providing standardized alignments and an import filter
13 (greengenes.ift) that imports the alignment and other standardized fields from 16S
14 records vetted weekly from GenBank.

15 To exemplify the utility of the Greengenes data assembly process, and to
16 interrogate the validity of prokaryotic candidate phyla, we aligned and chimera-checked
17 over 90,000 public 16S sequences. Taxonomic classifications were applied from the
18 major curators where available. Sequence data were imported from NCBI for complete
19 or nearly complete gene sequences (>1250 nt in length) deposited as of April 2, 2006.
20 Alignment of both archaeal and bacterial sequences was performed with the NAST
21 aligner (8) against a “Core Set” of templates selected from a phylogenetically broad
22 collection (16). The resulting multiple sequence alignment (MSA) was formatted so that
23 each sequence occupies a consistent 7,682 characters or 4,182 characters, the later

1 enabling compatibility with RDP v8.1(22) alignments. Both these formats were concise
2 enough for browsing in common MSA graphical interfaces such as ClustalX(28),
3 MEGA(18) and the platform independent Jalview(6) as well as ARB. Other standard
4 expansions such as the >20,000-character Ludwig alignment will be an alternate format
5 available in future releases to give maximum flexibility to researchers.

6 For high-throughput chimera screening of the aligned sequences, the program
7 Bellerophon (15) was used with two modifications. First, the algorithm was modified to
8 reduce the number of potential parents considered in the partial trees allowing run time to
9 scale linearly rather than logarithmically with the count of candidate sequences in a
10 collection. Secondly, a new metric was implemented that weights the likelihood of a
11 sequence being chimeric according to the similarity of the parent sequences. The more
12 distantly related the parent sequences are to each other relative to their divergence from
13 the candidate chimeric sequence, the greater the likelihood that the inferred chimera is
14 real. The metric, called the divergence ratio, uses the average sequence identity between
15 the two fragments of the candidate and their corresponding parent sequences as the
16 numerator, and the sequence identity between the parent sequences as the denominator.
17 All calculations were restricted to 1,287 conserved columns of aligned characters using a
18 300 base pair window on either side of the most likely break point. A divergence ratio >
19 1.1 and both fragment-to-parent similarities > 90% was required for classifying
20 sequences as putatively chimeric.

21 Taxonomy was linked to each record by various methods. NCBI and RDP
22 taxonomic nomenclature were extracted directly from the respective GenBank-formatted
23 records. The Pace and Ludwig annotations were exported from curated ARB databases.

1 The Hugenholtz taxonomy was also derived from a curated ARB database in which tree
2 topologies had been verified using RAxML-VI (27) for maximum likelihood inference.
3 The general time reversible model of evolution (GTR) was applied together with
4 optimization of substitution rates and site-specific rates according to a gamma
5 distribution. Different search algorithms were considered depending on the run-time of
6 the standard hill-climb (SHC) search method. If running in less than 8 hours, simulating
7 annealing (SA) was processed with the default starting temperature and a termination
8 time set to approximately 24 hours. If SA was not used and SHC terminated within 24
9 hours, SHC was used. Further, the rapid hill-climb (RHC) was used in all other cases
10 when its running time was less than 24 hours. If RHC did not terminate within the set
11 limit, the number of taxa was reduced. After 100 bootstrap replications, a consensus tree
12 was calculated using Consense (12) and imported into ARB. This database
13 (greengenes.arb) is available for download through Greengenes and is updated
14 periodically.

15 Of the 90,000 NCBI records analyzed, 54% were derived from uncultured
16 organisms, the majority of which have been deposited over the last five years (Figure 1).
17 Only three studies have submitted greater than one thousand full-length clones, however,
18 we expect the number of large 16S surveys to rise due to the availability and falling cost
19 of high throughput sequencing. Bellerophon detection of putative chimeras in 3% of the
20 sequences from uncultured organisms was not unexpected considering initial estimates
21 (2). Surprisingly, 0.2% of sequences derived from pure cultures were also determined to
22 be putative chimeras. Multiple distinct 16S rRNA genes have been encountered when
23 creating clone libraries from colonies assumed to be pure cultures prepared from

1 numerous third party sources (Colleen Cavanaugh, personal communication). There is a
2 possibility that isolated colonies contain symbiotic bacteria which increase PCR template
3 complexity enabling chimera formation. In addition, thousands of full-length, 16S-
4 annotated GenBank records only partially aligned using NAST. Future versions of
5 NAST could be altered to allow alignment extensions across regions of low template
6 similarity or to allow candidates to be aligned in sections using divergent templates. Both
7 of these options may allow a greater abundance of chimeric data to be imported into
8 Greengenes, but, perhaps, would capture novel phyla from the public repositories.
9 Alternately, manually aligned sequences from novel phyla can be offered from the user
10 community for recruitment to the Core Set advocating periodic re-evaluation of the
11 partially aligned set.

12 Discovery of chimeras in 16S data collections is crucial if the data set will be a
13 foundation for applied bioinformatics. They pose a fundamental problem when used as
14 templates with probe selection software, a growing concern with the recent increase in
15 16S microarray probe development (3, 8, 11). The 15 to 30 bases surrounding the
16 chimeric break-point can appear sufficiently unique from all other records in a database
17 causing a probe selection algorithm to justifiably identify this region as a target's
18 signature and suggest complementary probes to be synthesized. These probes would
19 appear highly valuable considering their minimal mis-hybridization potential, but, in fact,
20 they would rarely be useful since they target non-existent organisms. Chimera test results
21 from Greengenes allow greater control over input to probe selection software, will aid in
22 avoiding artificial T-RFLP pattern predictions from the ARB-compatible TRF-CUT (25)
23 and can increase the accuracy of sampling rarefaction curves (26).

1 The fraction of putative chimeras within the deposited sequences from an
2 individual study varies from none to over 20% (Figure 1) suggesting that chimera
3 screening is still not being uniformly applied by sequence generators. The problem is
4 exacerbated in sparsely populated candidate phyla. For instance, bacterial phyla, ‘SAM’,
5 and ‘5’ and the class GN4 (Proteobacteria) may require re-evaluation. Likewise the
6 genera *Tistrella*, *Caldotoga*, *Dehalobacterium* and *Desulfovermiculus* are currently
7 anchored by sequences with evidence of chimeric composition. Additional sequences
8 could lead to the empirical rejection of certain classifications or may aid in defining the
9 true breadth of sequence variation for these taxa.

10 Comparison of five different taxonomies uncovered a surprisingly large disparity
11 between expert curators. Loosely interpreting “phyla” to be any labeled grouping or
12 division immediately subordinate to the domains Archaea or Bacteria, the five curations
13 were compared in a Venn diagram (Figure 2). The main source of the disparity arises
14 from discordant naming of novel candidate phyla, or absence of names for candidate
15 phyla. For example Pace and Hugenholtz have independently named over a dozen
16 phylum-level lineages, many of which are the same lineages, and RDP has not named any
17 of these lineages. This is a consequence of the huge number of environmental sequences
18 in the public databases and frequent redundant naming of environmental lineages in the
19 literature. We hope that making multiple taxonomic classifications available through
20 Greengenes will aid in standardizing classification, particularly classification of
21 environmental lineages.

22 Greengenes is also a functional workbench to assist in analysis of user-generated
23 16S rRNA gene sequences. Batches of sequencing reads can be uploaded for quality-

1 based trimming and creation of multiple sequence alignments (9). Three types of non-
2 MSA similarity searches are also available; seed extension by BLAST (1), similarity
3 based on shared 7-mers by a tool entitled “Simrank”, or direct degenerative pattern match
4 for probe/primer evaluation. Results are displayed using user-preferred taxonomic
5 nomenclature and can be saved between sessions.

6 In summary, Greengenes offers annotated, chimera-checked, full-length 16S
7 rRNA gene sequences in standard alignment formats. The relational database links
8 taxonomies from multiple curators and multiple sequences from a single study. It was
9 revealed that incongruent taxonomic nomenclature exists among curators even at the
10 phylum-level. Bellerophon found putative chimeras in sequences derived from both
11 uncultured and isolated organisms. The data set can be compared to user-provided
12 sequences via web interface or can be imported directly into ARB for advanced analyses.
13 We anticipate that Greengenes will be valuable to researchers conducting environmental
14 surveys and in 16S rRNA microarray design.

15 In the immediate future, we plan to develop and implement a number of
16 community curation tools. This will allow the user community to actively participate in
17 improving the quality of the Greengenes database and ensure that time-consuming
18 manual improvements of sequence and sequence-associated data including taxonomic
19 corrections are propagated for the benefit of the whole community. Specifically five
20 curation tools are in development that should capture manual improvements: 1)
21 improvements in individual sequence alignments, 2) manual verification of putative
22 chimeras, 3) recruitment of novel lineages to the Core Set, 4) corrections in the
23 Greengenes description (the abbreviated description of the record usually taking the form

1 [habitat] clone [clone name] for environmental sequences), and 5) updating taxonomic
2 group names. One of the main challenges in the implementation of these tools will be to
3 ensure that only high quality manual edits are incorporated into Greengenes. For
4 example, for a suggested alignment alteration, the submitted sequence must a) match the
5 existing sequence, b) preserve the location of highly conserved positions in the 16S
6 rRNA gene and c) record the curator information as part of the update transaction. We
7 recognize a desire on the part of many users to contribute to a distributed curation effort
8 and we hope that Greengenes will become a resource to facilitate this aim.

9 **Acknowledgments**

10 We thank Kirk Harris and Norm Pace for sharing their ARB database and Richard
11 Phan for assistance with web graphics. The computational infrastructure was provided in
12 part by the Virtual Institute for Microbial Stress and Survival (<http://VIMSS.lbl.gov>)
13 supported by the U. S. Department of Energy, Office of Science, Office of Biological and
14 Environmental Research, Genomics:GTL Program and the Natural and Accelerated
15 Bioremediation Research Program through contract DE-AC02-05CH11231 between
16 Lawrence Berkeley National Laboratory and the U. S. Department of Energy. Web
17 application development was funded in part by the Department of Homeland Security
18 under grant number HSSCHQ04X00037.
19

References

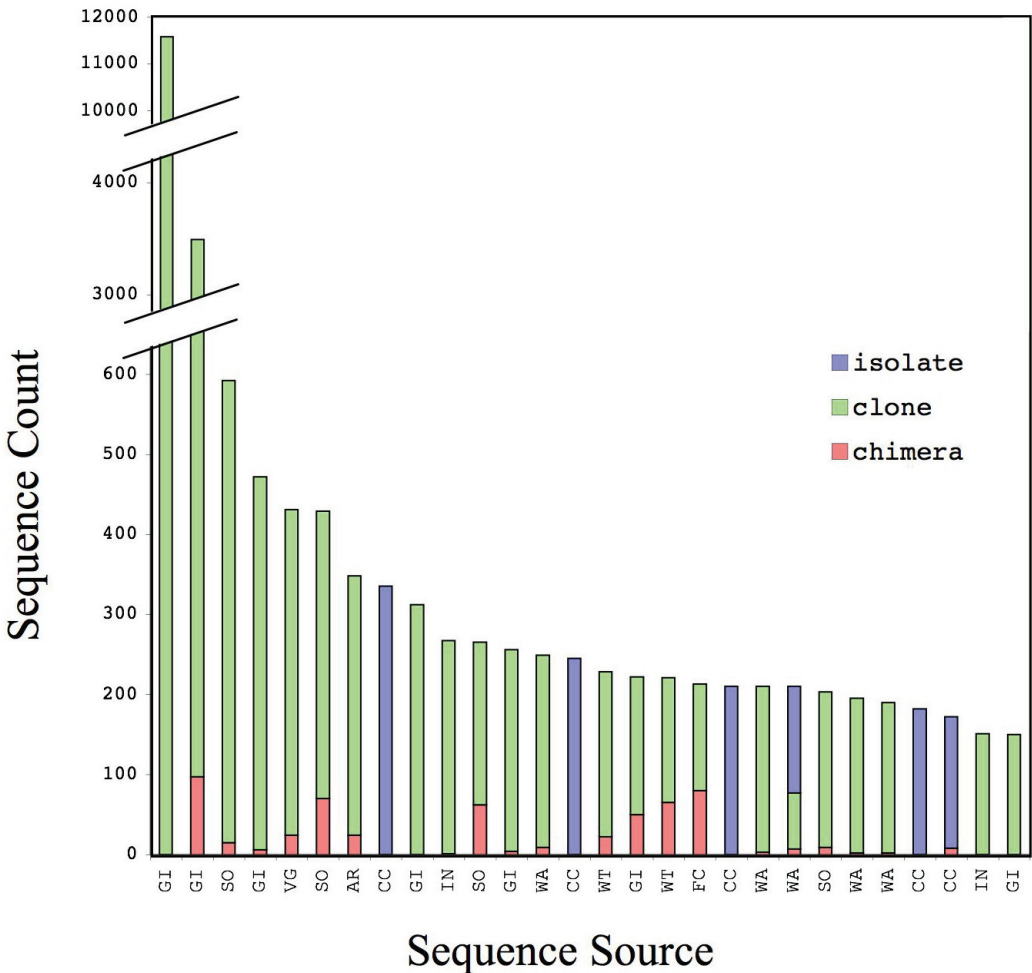
1. **Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman.** 1990. Basic local alignment search tool. *J Mol Biol* **215**:403-10.
2. **Ashelford, K. E., N. A. Chuzhanova, J. C. Fry, A. J. Jones, and A. J. Weightman.** 2005. At least 1 in 20 16S rRNA sequence records currently held in public repositories is estimated to contain substantial anomalies. *Appl Environ Microbiol* **71**:7724-36.
3. **Ashelford, K. E., A. J. Weightman, and J. C. Fry.** 2002. PRIMROSE: a computer program for generating and estimating the phylogenetic range of 16S rRNA oligonucleotide probes and primers in conjunction with the RDP-II database. *Nucleic Acids Res* **30**:3481-9.
4. **Brodie, E., S. Edwards, and N. Clipson.** 2002. Bacterial community dynamics across a floristic gradient in a temperate upland grassland ecosystem. *Microb Ecol* **44**:260-70.
5. **Castiglioni, B., E. Rizzi, A. Frosini, K. Sivonen, P. Rajaniemi, A. Rantala, M. A. Mugnai, S. Ventura, A. Wilmotte, C. Boutte, S. Grubisic, P. Balthasart, C. Consolandi, R. Bordoni, A. Mezzelani, C. Battaglia, and G. De Bellis.** 2004. Development of a universal microarray based on the ligation detection reaction and 16S rRNA gene polymorphism to target diversity of cyanobacteria. *Appl Environ Microbiol* **70**:7161-72.
6. **Clamp, M., J. Cuff, S. M. Searle, and G. J. Barton.** 2004. The Jalview Java alignment editor. *Bioinformatics* **20**:426-7.
7. **Cole, J. R., B. Chai, R. J. Farris, Q. Wang, S. A. Kulam, D. M. McGarrell, G. M. Garrity, and J. M. Tiedje.** 2005. The Ribosomal Database Project (RDP-II): sequences and tools for high-throughput rRNA analysis. *Nucleic Acids Res* **33**:D294-6.
8. **DeSantis, T. Z., I. Dubosarskiy, S. R. Murray, and G. L. Andersen.** 2003. Comprehensive aligned sequence construction for automated design of effective probes (CASCADE-P) using 16S rDNA. *Bioinformatics* **19**:1461-8.
9. **DeSantis, T. Z., P. Hugenholtz, K. Keller, E. L. Brodie, N. Larsen, Y. M. Piceno, R. Phan, and G. L. Andersen.** 2006. NAST: A multiple sequence alignment server for comparative analysis of 16S rRNA genes. *Nucleic Acids Res*:accepted.
10. **DeSantis, T. Z., C. E. Stone, S. R. Murray, J. P. Moberg, and G. L. Andersen.** 2005. Rapid quantification and taxonomic classification of environmental DNA from both prokaryotic and eukaryotic origins using a microarray. *FEMS Microbiol Lett* **245**:271-8.
11. **Emrich, S. J., M. Lowe, and A. L. Delcher.** 2003. PROBEmer: A web-based software tool for selecting optimal DNA oligos. *Nucleic Acids Res* **31**:3746-50.
12. **Felsenstein, J.** 1989. PHYLIP - Phylogeny Inference Package (Version 3.65). *Cladistics* **5**:164-166.
13. **Harris, J. K., S. T. Kelley, and N. R. Pace.** 2004. New perspective on uncultured bacterial phylogenetic division OP11. *Appl Environ Microbiol* **70**:845-9.
14. **Harris, J. K., S. T. Kelley, G. B. Spiegelman, and N. R. Pace.** 2003. The genetic core of the universal ancestor. *Genome Res* **13**:407-12.

- 1 15. **Huber, T., G. Faulkner, and P. Hugenholtz.** 2004. Bellerophon: a program to
2 detect chimeric sequences in multiple sequence alignments. *Bioinformatics*
3 **20**:2317-9.
- 4 16. **Hugenholtz, P.** 2002. Exploring prokaryotic diversity in the genomic era.
5 *Genome Biol* **3**:1-8.
- 6 17. **Kelly, J. J., S. Siripong, J. McCormack, L. R. Janus, H. Urakawa, S. El**
7 **Fantroussi, P. A. Noble, L. Sappelsa, B. E. Rittmann, and D. A. Stahl.** 2005.
8 DNA microarray detection of nitrifying bacterial 16S rRNA in wastewater
9 treatment plant samples. *Water Res* **39**:3229-38.
- 10 18. **Kumar, S., K. Tamura, and M. Nei.** 2004. MEGA3: Integrated software for
11 Molecular Evolutionary Genetics Analysis and sequence alignment. *Brief*
12 *Bioinform* **5**:150-63.
- 13 19. **Lane, D. J., A. P. Harrison, Jr., D. Stahl, B. Pace, S. J. Giovannoni, G. J.**
14 **Olsen, and N. R. Pace.** 1992. Evolutionary relationships among sulfur- and iron-
15 oxidizing eubacteria. *J Bacteriol* **174**:269-78.
- 16 20. **Lehner, A., A. Loy, T. Behr, H. Gaenge, W. Ludwig, M. Wagner, and K. H.**
17 **Schleifer.** 2005. Oligonucleotide microarray for identification of *Enterococcus*
18 species. *FEMS Microbiol Lett* **246**:133-42.
- 19 21. **Ludwig, W., O. Strunk, R. Westram, L. Richter, H. Meier, Yadhukumar, A.**
20 **Buchner, T. Lai, S. Steppi, G. Jobb, W. Forster, I. Brettske, S. Gerber, A. W.**
21 **Ginhart, O. Gross, S. Grumann, S. Hermann, R. Jost, A. Konig, T. Liss, R.**
22 **Lusmann, M. May, B. Nonhoff, B. Reichel, R. Strehlow, A. Stamatakis, N.**
23 **Stuckmann, A. Vilbig, M. Lenke, T. Ludwig, A. Bode, and K. H. Schleifer.**
24 2004. ARB: a software environment for sequence data. *Nucleic Acids Res*
25 **32**:1363-71.
- 26 22. **Maidak, B. L., J. R. Cole, T. G. Lilburn, C. T. Parker, Jr., P. R. Saxman, R.**
27 **J. Farris, G. M. Garrity, G. J. Olsen, T. M. Schmidt, and J. M. Tiedje.** 2001.
28 The RDP-II (Ribosomal Database Project). *Nucleic Acids Res* **29**:173-4.
- 29 23. **Pace, N. R.** 1997. A molecular view of microbial diversity and the biosphere.
30 *Science* **276**:734-40.
- 31 24. **Radosevich, J. L., W. J. Wilson, J. H. Shinn, T. Z. DeSantis, and G. L.**
32 **Andersen.** 2002. Development of a high-volume aerosol collection system for the
33 identification of air-borne micro-organisms. *Lett Appl Microbiol* **34**:162-7.
- 34 25. **Ricke, P., S. Kolb, and G. Braker.** 2005. Application of a newly developed
35 ARB software-integrated tool for in silico terminal restriction fragment length
36 polymorphism analysis reveals the dominance of a novel *pmoA* cluster in a forest
37 soil. *Appl Environ Microbiol* **71**:1671-3.
- 38 26. **Schloss, P. D., and J. Handelsman.** 2004. Status of the microbial census.
39 *Microbiol Mol Biol Rev* **68**:686-91.
- 40 27. **Stamatakis, A., T. Ludwig, and H. Meier.** 2005. RAXML-III: a fast program for
41 maximum likelihood-based inference of large phylogenetic trees. *Bioinformatics*
42 **21**:456-63.
- 43 28. **Thompson, J. D., T. J. Gibson, F. Plewniak, F. Jeanmougin, and D. G.**
44 **Higgins.** 1997. The CLUSTAL_X windows interface: flexible strategies for
45 multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res*
46 **25**:4876-82.

- 1 29. **Webster, G., C. J. Newberry, J. C. Fry, and A. J. Weightman.** 2003.
2 Assessment of bacterial community structure in the deep sub-seafloor biosphere
3 by 16S rDNA-based techniques: a cautionary tale. *J Microbiol Methods* **55**:155-
4 64.
- 5 30. **Wilson, K. H., W. J. Wilson, J. L. Radosevich, T. Z. DeSantis, V. S.**
6 **Viswanathan, T. A. Kuczmarski, and G. L. Andersen.** 2002. High-density
7 microarray of small-subunit ribosomal DNA probes. *Appl Environ Microbiol*
8 **68**:2535-41.
9
10

1
2
3
4
5
6
7

Figure 1. 16S rRNA gene sequencing projects producing over 200 full-length records. All projects were submitted to GenBank between October 2000 and February 2006. Sequences were generated from gastrointestinal (GI), soil (SO), vaginal (VG), aerosol (AR), culture collection (CC), insect (IN), water (WA), waste treatment (WT) or fecal (FC) sources as labeled on the x-axis and. Projects are ordered by sequence count.



8
9
10

Figure 2. Phyla-level nomenclature shared among independent curators represented as a five-way Venn diagram. Yellow spheres represent the 126 Phyla or Candidate Division names encountered in at least one of the five taxonomy systems (Pace, Hugenholtz, Ludwig, RDP, or NCBI). Numbers in parentheses are the count of Phyla or Candidate Division names recognized by an individual curator. Clusters of yellow spheres connected by more than one colored web symbolize names recognized by multiple curators. Image rendered by AutoFocus software (Aduna B.V., The Netherlands). A complete table of phyla-level nomenclature comparisons is available at: <http://greengenes.lbl.gov/TaxCompare>.

